

EUROPEAN PATENT OFFICE

Patent Abstracts of Japan

PUBLICATION NUMBER : 11282492
PUBLICATION DATE : 15-10-99

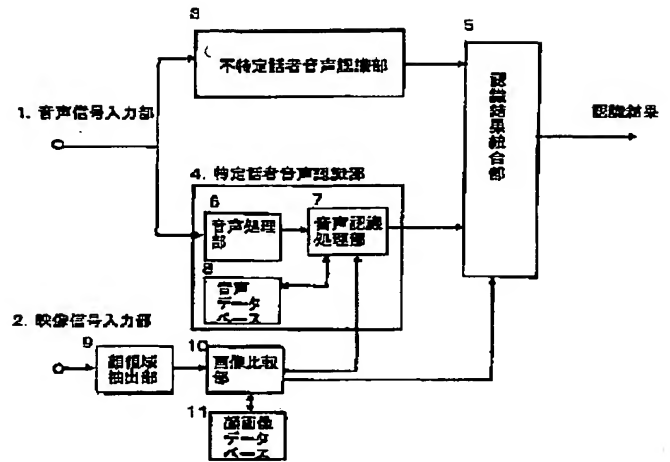
APPLICATION DATE : 26-03-98
APPLICATION NUMBER : 10079916

APPLICANT : MATSUSHITA ELECTRIC IND CO LTD;

INVENTOR : INOUE IKUO;

INT.CL. : G10L 3/00 G10L 3/00 G10L 3/00
G06T 1/00 G10L 5/06

TITLE : SPEECH RECOGNITION DEVICE,
SPEAKER DETECTOR, AND IMAGE
RECORDER



ABSTRACT : PROBLEM TO BE SOLVED: To achieve speech recognition of plural speakers with high reliability.

SOLUTION: This system is configured of a speech signal input part 1, a video signal input part 2, an unspecified speaker speech recognition part 3 for extracting a common feature from speeches of multi-speakers, making a standard pattern, and calculating a degree of similarity between the input speeches and a standard speech pattern, a specific speaker speech recognition part 4 for calculating a degree of similarity between the input speech and the speech of a pre-registered speaker, a face region extracting part 9 for extracting a face region from an input video, a face image database 11 for recording face image data of plural specific speakers and their identification numbers, an image comparison part 10 for outputting the degree of similarity with the image data inputted from the face region extracting part 9 and the face image database 11, and a recognition result integration part 5 for calculating an integrated degree of similarity from the outputs of the unspecified speaker speech recognition part 3, the specific speaker speech recognition part 4, and the image comparison part 10, and outputting the recognition result.

COPYRIGHT: (C)1999,JPO

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-282492

(43) 公開日 平成11年(1999)10月15日

(51) Int.Cl.⁶

G 1 0 L 3/00

識別記号

5 7 1

5 1 3

5 3 1

F I

G 1 0 L 3/00

5 7 1 C

5 1 3 Z

5 3 1 J

5 3 1 K

G 0 6 T 1/00

5/06

D

審査請求 未請求 請求項の数12 O L (全 11 頁) 最終頁に続く

(21) 出願番号

特願平10-79916

(22) 出願日

平成10年(1998) 3 月26日

(71) 出願人 000005821

松下電器産業株式会社

大阪府門真市大字門真1006番地

(72) 発明者 古 山 浩 志

大阪府門真市大字門真1006番地 松下電器
産業株式会社内

(72) 発明者 井 上 郁 夫

大阪府門真市大字門真1006番地 松下電器
産業株式会社内

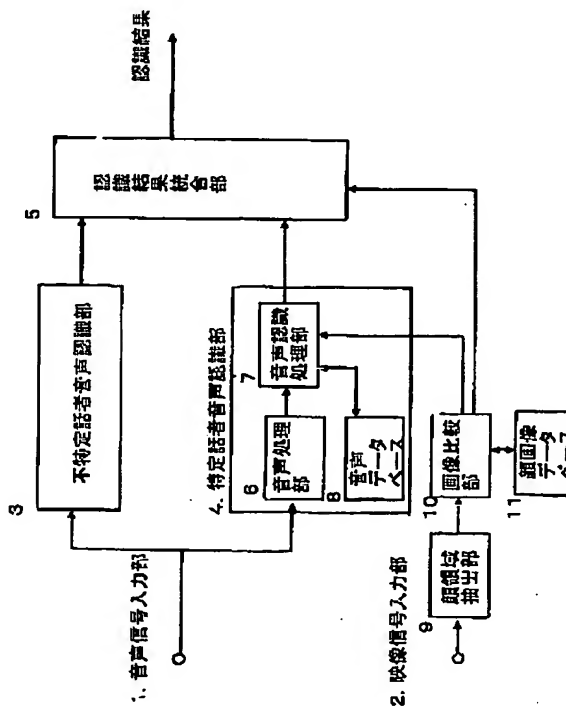
(74) 代理人 弁理士 蔵合 正博

(54) 【発明の名称】 音声認識装置、話者検出装置及び画像記録装置

(57) 【要約】

【課題】 複数の話者に対して、信頼性の高い音声認識を実現する。

【解決手段】 音声信号入力部1と、映像信号入力部2と、複数話者の音声から共通する特徴を抽出して標準パターンを作成し入力音声と標準音声パターンとの類似度を算出する不特定話者音声認識部3と、予め登録された話者の音声と入力音声との類似度を算出する特定話者音声認識部3と、入力映像から顔領域を抽出する顔領域抽出部9と、複数特定話者の顔画像データと話者識別番号とを記録する顔画像データベース11と、顔領域抽出部9と顔画像データベース11から入力する画像データとの類似度を出力する画像比較部10と、不特定話者音声認識部3、特定話者音声認識部4、及び画像比較部10の出力から統合的類似度を算出し認識結果を出力する認識結果統合部5とで構成される。



【特許請求の範囲】

【請求項1】 話者の特徴的外観の画像を含む映像データを入力する映像入力手段と、話者の音声データを入力する音声入力手段と、

複数の特定の話者の音声データを、それを特定できる話者識別情報とともに登録し、登録された音声データと入力音声データとの類似度を算出して音声認識を行う特定話者音声認識手段と、

不特定多数の話者の音声データから共通する特徴を抽出して標準パターンを作成、登録して、音声標準パターンと入力音声データとの類似度を算出して音声認識を行う不特定話者音声認識手段と、

複数の話者の特徴的外観の画像を、その話者を特定できる話者識別情報とともに登録し、登録された画像データと前記映像データに含まれる話者の画像データとの類似度を算出する画像認識手段と、

特定話者音声認識手段の出力と不特定話者音声認識手段からの出力と画像認識手段からの出力とを統合して、音声認識結果として単語等を出力する認識結果統合手段とを備えた音声認識装置。

【請求項2】 話者の特徴的外観を含む映像データを入力する映像入力手段と、

話者の音声データを入力する音声入力手段と、

複数の話者の特徴的外観の画像を、その話者を特定できる話者識別情報とともに登録し、登録された画像データと前記映像データに含まれる話者の画像データとの類似度を算出する画像認識手段と、

複数の特定の話者の音声データを、それを特定できる情報とともに登録し、前記画像認識手段で算出された類似度をもとに登録された音声データを絞り込んだ後、その絞り込まれた音声データと入力音声データとの類似度を算出して音声認識を行う特定話者音声認識手段と、不特定多数の話者の音声データから共通する特徴を抽出して標準パターンを作成、登録して、音声標準パターンと入力音声データとの類似度を算出して音声認識を行う不特定話者音声認識手段と、

特定話者音声認識手段の出力と不特定話者音声認識手段からの出力とを統合して、音声認識結果として単語等を出力する認識結果統合手段とを備えた音声認識装置。

【請求項3】 特定話者音声認識手段では、画像認識手段から出力される類似度が閾値を超えた話者の音声データのみに対して、入力音声データとの類似度を算出する請求項1に記載の音声認識装置。

【請求項4】 特定話者音声認識手段では、画像認識手段から出力される類似度のうち最大となる話者の音声データのみに対して、入力音声データとの類似度を算出する請求項1に記載の音声認識装置。

【請求項5】 認識結果統合手段が、画像認識手段から出力される、話者識別情報に対応する話者の画像データと入力映像に含まれる画像データとの類似度を R_i 、特

定話者音声認識手段から出力される、話者識別情報に対応する話者の音声データ j に対する入力音声データとの類似度を $R'_{i,j}$ 、不特定話者音声認識手段から出力される、入力音声データと音声データ j との類似度を R''_j とするときに、その R_i と $R'_{i,j}$ と R''_j とを用いて最適な音声データを出力することを特徴とする請求項1、3、4のいずれかに記載の音声認識装置。

【請求項6】 顔などの唇を含む外観を話者の特徴的外観とし、入力映像から話者の唇の動きを検出する唇動き検出手段を備え、特定話者認識手段では、単位時間あたりの唇の動き量が設定された閾値よりも大きい入力に対してのみ、入力音声データと登録音声データの類似度を算出することを特徴とする請求項1、3～5のいずれかに記載の音声認識装置。

【請求項7】 話者の顔を含む映像データを入力する映像入力部と、話者の音声データを入力する音声入力部と、入力映像から話者の唇の動きを検出する唇動き検出手段と、入力音声から音声レベルを検出する音声レベル検出手段とを備え、単位時間あたりの唇の動き量と音声レベルが共に設定された閾値を超えている時には、入力映像中に話者の映像が含まれていることを示す話者検出信号を出力する話者検出装置。

【請求項8】 話者の顔を含む映像データを入力する映像入力部と、話者の音声データを入力する音声入力部と、入力映像から話者の唇の動きを検出する唇動き検出手段と、入力音声から音声レベルを検出する音声レベル検出手段とを備え、単位時間あたりの唇の動き量と音声レベルが共に設定された閾値を超えている時には、入力映像中に話者の映像が含まれていることを示す話者検出信号を出力する話者検出装置を具備し、特定話者認識手段では、話者検出信号が設定された閾値以上となる入力に対してのみ、入力音声データと登録音声データの類似度を算出することを特徴とする請求項1、3～6のいずれかに記載の音声認識装置。

【請求項9】 音声信号入力部と映像信号入力部は、それぞれ映像表示装置の音声信号出力部と映像出力部に接続され、前記映像表示装置の表示対象である出演者を特定できる出演者識別情報を含む、出演者情報を入力する出演者情報入力手段と、出演者情報を記録する出演者情報記録手段と、出演者情報から現在、表示されている出演者を特定し、登録された画像データの中から特定された出演者の画像データを検索する画像検索手段を備え、画像認識部では、検索された出演者の画像と入力した映像信号に含まれる話者の画像との類似度を算出することを特徴とする請求項1、3～6、8のいずれかに記載の音声認識装置。

【請求項10】 話者の特徴的外観を含む映像データを入力する映像入力部と、複数の話者の特徴的外観の画像を、それを特定できる話者識別情報と共に登録する画像

データベースと、登録された画像データと入力映像に含まれる話者の画像データとの類似度を算出する画像認識手段とを備え、

入力した映像に含まれる画像と登録された画像データの類似度を算出し、登録されているすべての画像データに対する類似度が予め設定された閾値以下の場合、未登録の話者として新たな話者識別情報とともに画像データベースに記録する画像記録装置。

【請求項11】 話者の特徴的外観を含む映像データを入力する映像入力部と、複数の話者の特徴的外観の画像を、それを特定できる話者識別情報と共に登録する画像データベースと、登録された画像データと入力映像に含まれる話者の画像データとの類似度を算出する画像認識手段とを備え、入力した映像に含まれる画像と登録された画像データの類似度を算出し、登録されているすべての画像データに対する類似度が予め設定された閾値以下の場合、未登録の話者として新たな話者識別情報とともに画像データベースに記録する画像記録装置を具備し、未登録の話者を自動的に登録することが可能な請求項1、3～6、8、9のいずれかに記載の音声認識装置。

【請求項12】 入力した映像に含まれる画像と登録された画像データの類似度を算出し、出力された類似度が予め設定した閾値S1以上となる登録話者のすべての音声データに対して、特定話者音声認識手段から出力する入力音声データと登録音声データとの間の類似度が予め設定された閾値S2以下であり、かつ、不特定話者音声認識手段から出力される候補単語等の類似度が予め設定された閾値S3以上である場合に、該当の話者の未登録音声データとして、それを特定できる話者識別情報とともに入力音声データを記録するための記録手段を有する請求項11に記載の音声認識装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、映像信号と音声信号を用いて音声認識を行う音声認識装置に関する。

【0002】

【従来の技術】音声認識方式には、特定話者音声認識方式と不特定話者音声認識方式とがある。特定話者の音声のみを認識する特定話者音声認識方式は、話者の音声を予め登録し、入力音声と登録音声の類似度を算出し、入力音声の認識を行うものである。従って、登録された話者の音声入力に対しては、高い認識率を得ることができるが、話者の音声を登録する作業が必要となる。一方、不特定話者認識方式では、不特定多数の音声から共通する特徴を抽出して標準パターンを作成し、入力音声と音声標準パターンとの類似度を算出し、入力音声の認識を行うものである。従って、話者の音声を登録する煩雑さはないが、特定話者音声認識方式と比較すると認識率は低くなる。

【0003】このような、特定話者音声認識方式と不特定話者音声認識方式における、それぞれの短所を補うため、特定話者音声認識方式と不特定話者音声認識方式を併用する音声認識装置が考えられている（特開昭63-32596号公報）。

【0004】この装置は図6に示すように、音声信号を入力する音声信号入力部1、不特定話者音声認識方式により音声認識を行う不特定話者認識部3、特定話者音声認識方式により音声認識を行う特定話者音声認識部4、不特定話者認識部3と特定話者認識部4でそれぞれ求めた認識結果を入力して、類似度の大きな方の認識結果を出力する認識結果統合部5を備えている。また、認識結果が正解と判断されたときには、入力音声の特徴データを特定話者認識用として特定話者認識部4に登録する。

【0005】このように、従来の音声認識装置では特定話者音声認識方式と不特定話者音声認識方式を併用することにより、音声認識の認識率を高め、また、特定の話者の音声データを自動的に登録することが可能となっている。

【0006】

【発明が解決しようとする課題】音声認識装置の用途として、例えばパーソナルコンピューター、TVやVTR等、家庭内にある電気製品の機器制御のための入力装置としての利用が考えられるが、家庭内で利用する場合には、ある特定の人物が発する音声に対してのみ高い認識率を有するのでは不十分であり、同居している家族など、複数の人物から発せられる、それぞれの音声入力に対しても高い認識率を維持する必要がある。

【0007】本発明は、このような要求にこたえるものであり、複数の話者に対しても高い認識率を実現することができる音声認識方式を提供することを目的としている。

【0008】

【課題を解決するための手段】そこで、本発明の音声認識装置では、話者の特徴的外観（顔など）を含む映像データを入力する映像入力手段と、話者の音声データを入力する音声入力手段と、不特定話者音声認識方式により音声認識を行う不特定話者音声手段と、認識を行う話者を含む複数の話者の音声データを蓄積する音声データベースと、特定話者音声認識方式により音声認識を行う特定話者音声認識手段と、入力する映像から話者の顔領域を抽出する顔領域抽出手段と、認識を行う話者を含む複数の話者の顔画像データを蓄積する顔画像データベースと、顔領域抽出手段から出力される顔画像と顔画像データベースに蓄積された顔画像とを比較して、類似度を出力する画像比較手段と、不特定話者音声認識手段と特定話者音声認識手段からそれぞれ出力される認識候補音声と入力音声との間の類似度と画像比較手段から出力される類似度を統合して、最終的な音声認識結果として出力する認識結果統合手段とを備え、不特定話者音声認

識手段から出力される入力音声と認識候補音声の類似度と、画像比較手段から出力される顔画像データベースに登録された話者の顔画像と入力映像に含まれる顔画像の類似度と、特定話者認識手段から出力される音声データベースに登録された話者の音声データと入力音声の類似度を組み合わせて、総合的な類似度から認識結果を出力するようにしている。

【0009】また、音声データベース、顔画像データベースに複数の話者のデータが登録されている場合には、それぞれの話者に対する顔画像と音声の類似度から、総合的な類似度を算出し認識結果を出力する。

【0010】従って、登録された話者の顔画像と入力映像に含まれる話者の顔画像の類似度が小さいときには、不特定話者音声認識手段からの出力が認識結果に大きく寄与し、登録された話者の顔画像と入力した話者の顔画像の類似度が大きいときには、その中でも最も類似度の大きな話者に対する、特定話者音声認識手段からの出力が認識結果に大きく寄与するため、複数の特定話者、あるいは不特定の話者から発せられる音声の入力に対して、より信頼性の高い音声認識が可能となる。

【0011】

【発明の実施の形態】本発明の請求項1に記載の発明は、音声認識装置に、話者の特徴的外観の画像を含む映像データを入力する映像入力手段と、話者の音声データを入力する音声入力手段と、複数の特定の話者の音声データを、それを特定できる話者識別情報とともに登録し、登録された音声データと入力音声データとの類似度を算出して音声認識を行う特定話者音声認識手段と、不特定多数の話者の音声データから共通する特徴を抽出して標準パターンを作成、登録して、音声標準パターンと入力音声データとの類似度を算出して音声認識を行う不特定話者音声認識手段と、複数の話者の特徴的外観の画像を、その話者を特定できる話者識別情報とともに登録し、登録された画像データと前記映像データに含まれる話者の画像データとの類似度を算出する画像認識手段と、特定話者音声認識手段の出力と不特定話者音声認識手段からの出力と画像認識手段からの出力とを統合して、音声認識結果として単語等を出力する認識結果統合手段とを備えたものであり、登録された話者の顔画像と入力映像に含まれる話者の顔画像の類似度の大小によって出力元の音声認識手段を変えることにより、複数の特定話者、あるいは不特定の話者から発せられる音声の入力に対して、より信頼性の高い音声認識が可能になるという作用を有する。

【0012】本発明の請求項2に記載の発明は、音声認識装置に、話者の特徴的外観を含む映像データを入力する映像入力手段と、話者の音声データを入力する音声入力手段と、複数の話者の特徴的外観の画像を、その話者を特定できる話者識別情報とともに登録し、登録された画像データと前記映像データに含まれる話者の画像デー

タとの類似度を算出する画像認識手段と、複数の特定の話者の音声データを、それを特定できる情報とともに登録し、前記画像認識手段で算出された類似度をもとに登録された音声データを絞り込んだ後、その絞り込まれた音声データと入力音声データとの類似度を算出して音声認識を行う特定話者音声認識手段と、不特定多数の話者の音声データから共通する特徴を抽出して標準パターンを作成、登録して、音声標準パターンと入力音声データとの類似度を算出して音声認識を行う不特定話者音声認識手段と、特定話者音声認識手段の出力と不特定話者音声認識手段からの出力とを統合して、音声認識結果として単語等を出力する認識結果統合手段とを備えたものであり、顔などを含む映像から話者の顔画像を抽出して、登録された話者の顔画像データベースと照合し、類似度を算出して、特定話者音声認識部、不特定話者音声認識部から出力する音声の類似度との総合的な類似度を算出して認識結果を出力することにより、複数の特定話者の入力に対して、信頼性の高い音声認識を行うことが可能となるという作用を有する。

【0013】本発明の請求項3に記載の発明は、請求項1記載の音声認識装置において、特定話者音声認識手段では、画像認識手段から出力される類似度が閾値を超えた話者の音声データのみに対して、入力音声データとの類似度を算出するようにしたものである。

【0014】本発明の請求項4に記載の発明は、請求項1記載の音声認識装置において、特定話者音声認識手段では、画像認識手段から出力される類似度のうち最大となる話者の音声データのみに対して、入力音声データとの類似度を算出するようにしたものである。

【0015】本発明の請求項5に記載の発明は、請求項1、3、4のいずれかに記載の音声認識装置において、認識結果統合手段が、画像認識手段から出力される、話者識別情報に対応する話者の画像データと入力映像に含まれる画像データとの類似度を R_i 、特定話者音声認識手段から出力される、話者識別情報に対応する話者の音声データ j に対する入力音声データとの類似度を $R'_{i,j}$ 、不特定話者音声認識手段から出力される、入力音声データと音声データ j との類似度を R''_j とするときに、その R_i と $R'_{i,j}$ と R''_j とを用いて最適な音声データを出力するようにしたものである。

【0016】本発明の請求項6に記載の発明は、請求項1、3～5のいずれかに記載の音声認識装置において、顔などの唇を含む外観を話者の特徴的外観とし、入力映像から話者の唇の動きを検出する唇動き検出手段を備え、特定話者認識手段では、単位時間あたりの唇の動き量が設定された閾値よりも大きい入力に対してのみ、入力音声データと登録音声データの類似度を算出するようにしたものである。

【0017】本発明の請求項7に記載の発明は、話者検出装置として、話者の顔を含む映像データを入力する映

像入力部と、話者の音声データを入力する音声入力部と、入力映像から話者の唇の動きを検出する唇動き検出手段と、入力音声から音声レベルを検出する音声レベル検出手段とを備えたものであり、単位時間あたりの唇の動き量と音声レベルが共に設定された閾値を超えている時には、入力された映像中に話者の映像が含まれていることを示す話者検出信号を出力するという作用を有する。

【0018】本発明の請求項8に記載の発明は、請求項1、3～6のいずれかに記載の音声認識装置において、請求項7に記載の話者検出装置と同様の機能を有する話者検出手段を具備し、特定話者認識手段では、話者検出信号が設定された閾値以上となる入力に対してのみ、入力音声データと登録音声データの類似度を算出するようにしたものである。

【0019】本発明の請求項9に記載の発明は、請求項1、3～6、8のいずれかに記載の音声認識装置において、音声信号入力部と映像信号入力部は、それぞれ映像表示装置の音声信号出力部と映像出力部に接続され、前記映像表示装置の表示対象である出演者を特定できる出演者識別情報を含む、出演者情報を入力する出演者情報入力手段と、出演者情報を記録する出演者情報記録手段と、出演者情報から現在、表示されている出演者を特定し、登録された画像データの中から特定された出演者の画像データを検索する画像検索手段を備えたものであり、画像認識部では、検索された出演者の画像と入力した映像信号に含まれる話者の画像との類似度を算出し、類似度に応じた音声を出力するという作用を有する。

【0020】本発明の請求項10に記載の発明は、話者の特徴的外観を含む映像データを入力する映像入力部と、複数の話者の特徴的外観の画像を、それを特定できる話者識別情報と共に登録する画像データベースと、登録された画像データと入力映像に含まれる話者の画像データとの類似度を算出する画像認識手段とを備えたものであり、入力した映像に含まれる画像と登録された画像データの類似度を算出し、登録されているすべての画像データに対する類似度が予め設定された閾値以下の場合には、未登録の話者として新たな話者識別情報とともに画像データベースに記録することによりデータベースの内容をより充実させるという作用を有する。

【0021】本発明の請求項11に記載の発明は、請求項1、3～6、8、9のいずれかに記載の音声認識装置において、請求項10に記載の顔画像装置と同様の機能を有する画像記録手段を具備したものであり、未登録の話者を自動的に登録することによりデータベースの内容をより充実させるという作用を有する。

【0022】本発明の請求項12に記載の発明は、請求項11記載の音声認識装置において、入力した映像に含まれる画像と登録された画像データの類似度を算出し、出力された類似度が予め設定した閾値S1以上となる登

録話者のすべての音声データに対して、特定話者音声認識手段から出力する入力音声データと登録音声データとの間の類似度が予め設定された閾値S2以下であり、かつ、不特定話者音声認識手段から出力される候補単語等の類似度が予め設定された閾値S3以上である場合に、該当の話者の未登録音声データとして、それを特定できる話者識別情報とともに入力音声データを記録するための記録手段を備えたものである。

【0023】以下、本発明の実施の形態について図1から図5を用いて説明する。なお、本発明はこれら実施の形態に何等限定されるものではなく、その要旨を逸脱しない範囲において種々なる態様で実施し得る。

【0024】(実施の形態1) 第一の実施の形態の音声認識装置は、図1に示すように、音声を入力する音声信号入力部1、映像を入力する映像信号入力部2、複数の不特定話者の音声から共通する特徴を抽出して標準パターンを作成し、入力音声と音声標準パターンとの間の類似度を算出して出力する不特定話者音声認識部3、予め登録された話者の音声と入力音声の類似度を算出して出力する特定話者音声認識部4、入力映像から話者の顔領域を抽出する顔領域抽出部9、複数の特定話者の顔画像データを話者の名前、もしくは話者識別コードとともに記録する顔画像データベース11、顔領域抽出部9と顔画像データベース11から入力する画像データを比較し、類似度を出力する画像比較部10、不特定話者認識部3、特定話者認識部4、画像比較部10から出力されるそれぞれの類似度から、類似度がもっとも大きな類似度に対応する単語等を認識結果として出力する認識結果統合部5により構成される。

【0025】また、特定話者音声認識部4は、入力した音声スペクトラム分析等により音声の特徴量を抽出する音声処理部6、隠れマルコフモデル等により入力音声と登録音声データとの間の特徴量の類似度を算出して出力する音声認識処理部7、複数の特定話者の音声データを話者の名前、もしくは話者識別コードとともに記録する音声データベース8により構成される。

【0026】この装置では、音声入力部1に入力した話者の音声は不特定話者音声認識部3と特定話者音声認識部4内の音声処理部6に入力する。音声処理部6の出力は、音声認識処理部7に入力し、音声認識処理部7の出力は、特定話者音声認識部4の出力として出力される。また、音声認識処理部7と音声データベース8は互いに接続されている。

【0027】映像信号入力部2に入力した話者の顔を含む映像は、顔領域抽出部9に入力し、顔領域抽出部9の出力は画像比較部10に入力し、画像比較部10の出力は音声認識処理部7に接続される。また、画像比較部10と画像データベース11は互いに接続されている。

【0028】不特定話者音声認識部3、特定話者音声認識部4、画像比較部10の出力は、認識結果統合部5に

入力し、認識結果統合部5からは認識結果が出力される。なお、映像信号中に含まれる人物の顔領域部分を抽出する手法は公知であり、例えば、第2回画像センシングシンポジウム講演集、A-1、pp. 1~6、「色情報とGAを用いた顔画像抽出と個人照合の応用」などに示されている。また、二つの顔画像データを比較する手法も、例えば、電子情報通信学会論文誌、D-2、Vol. 1., J76-D-2, No. 6, pp. 1132~1139、「モザイクとニューラルネットを用いた顔画像の認識」などに示されている。また、こうした技術を用いて、顔画像を個人照合に利用することは特願平8-170866号公報、或いは特願平8-86171号公報などに示されている。この実施の形態の装置においても、これらの技術を用いることが可能であるが、それだけに限定されるものではない。

【0029】また、音声認識処理部7における類似度の算出には、隠れマルコフモデル以外にニューラルネットワーク等、一般に用いられている他の手法を用いてもよい。顔画像比較部10では、登録された複数の話者に対して、i番目の話者の顔画像データと顔領域抽出部9から出力する顔画像データとの間の類似度 R_i を算出し出力する。特定話者音声認識部4では、登録された複数の話者に対してi番目の話者の音声データjと、入力音声との類似度を $R_{i,j}$ を算出して出力する。不特定話者認識部3では、音声データjに対して複数の不特定話者の音声から共通する特徴を抽出して作成した標準パターンと入力音声との間の類似度 $R'_{i,j}$ を算出し、出力する。認識結果統合部5では、あらかじめ設定した係数を α とするときに、登録された話者iのすべての音声データjに対して、 $\alpha \cdot R_i \cdot R'_{i,j}$ と $R'_{i,j}$ を算出し、その値が最大となる音声データjに対応する単語等を認識結果として出力する。

【0030】ここで、係数 α は、音声データベース8に登録した音声データとは別に、音声データベース8に登録された話者を含む複数の話者の音声データを用意し、この音声データの入力に対して認識率が最大となるようあらかじめ設定するものである。このように特定話者音声認識の方式において、登録された話者の顔画像データと入力映像に含まれる話者の顔画像を比較して、話者を特定することにより、信頼性の高い音声認識が可能となる。

【0031】（実施の形態2）この音声認識装置は図2に示すように、顔領域抽出部9から出力される話者の顔領域映像を入力し、前述の顔領域抽出部9と同じ手法により話者の口唇部分を抽出した映像を出力する唇領域抽出部12、唇領域抽出部12から出力される話者の口唇領域映像を入力し、話者の唇の動きを検出する唇動き検出部13と、実施の形態1と同じく音声処理部6、音声認識処理部7、音声データベース8から構成される特定話者音声認識部7、音声入力信号部1、映像信号入力部

2、不特定話者音声認識部3、顔領域抽出部9、画像比較部10、認識結果統合部5を備えている。

【0032】この装置では、音声入力部1に入力した話者の音声は不特定話者認識部3と特定話者認識部4内の音声処理部6に入力する。音声処理部6の出力は、音声認識処理部7に入力し、音声認識処理部7の出力は、特定話者音声認識部4の出力として出力される。また、音声認識処理部7と音声データベース8は互いに接続されている。映像信号入力部2に入力した話者の顔を含む映像は、顔領域抽出部9に入力し、顔領域抽出部9の出力は画像比較部10と唇領域抽出部12に入力し、画像比較部10の出力は音声認識処理部7に接続される。また、画像比較部10と画像データベース11とは互いに接続されている。

【0033】唇領域抽出部12の出力は唇動き検出部13に入力し、唇動き検出部13の出力は、音声認識処理部7に入力する。不特定話者音声認識部3、特定話者音声認識部4、画像比較部10の出力は、認識結果統合部5に入力し、認識結果統合部5からは認識結果が出力される。唇動き検出部13では、唇領域抽出部12で抽出した口唇領域の映像から、唇上のある着目点の動きベクトルを検出し、単位時間の唇の動きベクトルの平均値が設定した閾値よりも大きいときには

$$K=1$$

閾値より小さいときには

$$K=0$$

を出力する。認識結果統合部5では、登録されたすべての話者iの音声データjに対して、

$$\alpha \cdot K \cdot R_i \cdot R'_{i,j} \text{ と } R'_{i,j}$$

を算出し、その値が最大となる音声データjに対応する単語等を認識結果として出力する。

【0034】このように特定話者音声認識の方式において、入力映像に含まれる話者の唇の動きを検出することにより、入力映像に複数の話者の顔が含まれているときにも、信頼性の高い音声認識が可能となる。

【0035】（実施の形態3）この装置は図3に示すように、入力した音声のレベルが閾値を超えているかを検出する音声レベル検出部14と、実施の形態2と同じく、音声信号入力部1、音声処理部6、音声認識処理部7、音声データベース8から構成される特定話者音声認識部7、音声入力信号部1、映像信号入力部2、不特定話者音声認識部3、顔領域抽出部9、画像比較部10、認識結果統合部5、唇領域抽出部12、唇動き検出部13を備えている。

【0036】この装置では、音声入力部1に入力した話者の音声は不特定話者認識部3と特定話者認識部4内の音声処理部6に入力する。音声処理部6の出力は、音声認識処理部7に入力し、音声認識処理部7の出力は、特定話者音声認識部4の出力として出力される。また、音声認識処理部7と音声データベース8は互いに接続され

ている。映像信号入力部2に入力した話者の顔を含む映像は、顔領域抽出部9に入力し、顔領域抽出部9の出力は画像比較部10と唇領域抽出部12に入力し、画像比較部10の出力は音声認識処理部7に接続される。また、画像比較部10と画像データベース11は互いに接続されている。唇領域抽出部12の出力は唇動き検出部13に入力し、唇動き検出部13の出力は、音声認識処理部7に入力する。また、音声入力部1は音声レベル検出部14にも接続し、音声レベル検出部14の出力は音声認識処理部7と接続されている。

【0037】不特定話者音声認識部3、特定話者音声認識部4、画像比較部10の出力は、認識結果統合部5に入力し、認識結果統合部5からは認識結果が出力される。また、唇動き検出部13では、唇領域抽出部12で抽出した口唇領域の映像から、唇の動きを検出し、単位時間の唇の動きの平均値が設定した閾値よりも大きいときには

$K=1$

閾値よりも小さいときには

$K=0$

を出力する。音声レベル検出部14では、単位時間の音声レベルの平均値が設定した閾値よりも大きいときには

$L=1$

閾値よりも小さいときには

$L=0$

を出力する。認識結果統合部5では、登録されたすべての話者 i の音声データ j に対して、

$\alpha \cdot K \cdot L \cdot R_i \cdot R'_{i,j}$ と R''_j

を算出し、その値が最大となる音声データ j に対応する単語等を認識結果として出力する。

【0038】このように特定話者音声認識の方式において、入力映像に含まれる話者の唇の動きを検出することにより、入力映像に複数の話者の顔が含まれているときにも、より信頼性の高い音声認識が可能となる。

【0039】(実施の形態4) この装置は図4に示すように、TV番組の出演者の名前のデータを含む番組表を入力する番組表入力部15、番組表を記録する番組表記録部16、番組表データと現在の時刻を比較し、現在、放送されているTV番組の出演者を特定して出演者の名前を出力する出演者名検出部17、出演者名検出部17から出力する出演者の名前から顔画像データベース11を検索して、検索した顔画像を出力させる画像検索部18と、実施の形態1と同じく音声処理部6と音声データベース8と音声認識処理部7から構成される特定話者認識部4と、音声信号入力部1、映像信号入力部2と、不特定話者音声認識部3と、顔領域抽出部9と、画像比較部10と、認識結果統合部5を備えている。

【0040】この装置では、番組表入力部15に入力した番組表データは、番組表記録部16に入力し、記録される。出演者名検出部17は、番組表記録部16と画像

検索部18に接続され、画像検索部19と顔画像データベース11とは互いに接続されている。また、音声入力部1に入力した話者の音声は不特定話者認識部3と特定話者認識部4内の音声処理部6に入力する。音声処理部6の出力は、音声認識処理部7に入力し、音声認識処理部7の出力は、特定話者音声認識部4の出力として出力される。また、音声認識処理部7と音声データベース8は互いに接続されている。映像信号入力部2に入力した話者の顔を含む映像は、顔領域抽出部9に入力し、顔領域抽出部9の出力は画像比較部10に入力し、画像比較部10の出力は音声認識処理部7に接続される。また、画像比較部10と画像データベース11とは互いに接続されている。

【0041】不特定話者音声認識部3、特定話者音声認識部4、画像比較部10の出力は、認識結果統合部5に入力し、認識結果統合部5からは認識結果が出力される。また、音声信号入力部1と映像信号入力部2は、それぞれTV受信機の音声信号出力端子と映像出力端子に接続する。出演者名検出部17は、番組表記録部16に記録されているTV番組表データと現在の時刻から、現在放送されているTV番組の出演者を特定し、特定した出演者の名前データを出力する。画像検索部18は、時刻比較部17から出力する出演者の名前データをもとに画像検索データベース11から、出演者の顔画像を検索して、画像比較部12へと出力させる。顔画像比較部10では、TV番組の出演者と特定された複数の話者に対して、 k 番目の話者の顔画像データと入力映像に含まれる顔領域抽出部9から出力する顔画像データとの間の類似度 R_k を算出し出力する。特定話者音声認識部4では、出演者と特定された複数の話者に対して k 番目の話者の音声データ j と、入力音声との類似度を $R'_{k,j}$ を算出して出力する。不特定話者認識部3では、音声データ j に対して複数の不特定話者の音声から共通する特徴を抽出して作成した標準パターンと入力音声との間の類似度 R''_j を算出し、出力する。音声認識認識統合部では、あらかじめ設定した係数を α とするときに、出演者として特定されたすべての話者 k の音声データ j に対して、 $\alpha \cdot R_k \cdot R'_{k,j}$ と R''_j を算出し、その値が最大となる音声データ j に対応する単語等を認識結果として出力する。

【0042】このように特定話者音声認識の方式において、番組表データを本にTV番組の出演者を特定し、特定された出演者の顔画像データとTV番組映像に含まれる話者の顔画像を比較して、話者を特定することにより、より信頼性の高い音声認識が可能となる。なお、本実施の形態では映像表示装置としてTV受信機を例にあげて説明したが、VTR・ビデオ等の映像を表示するものであれば何でも構わない。

【0043】(実施の形態5) この装置は図5に示すように、複数の特定話者の音声データを話者の名前、もし

くは話者識別コードとともに記録し、かつ新規に音声データを追加記録する機能を有する音声データベース8、複数の特定話者の顔画像データを話者の名前、もしくは話者識別コードとともに記録し、かつ新規に顔画像データを追加記録する機能を有する顔画像データベース11、音声データベース8、顔画像データベース11にデータを追加記録するための制御を行う記録制御部19を備え、第一の実施の形態と同じく、音声処理部6、音声データベース8、音声認識処理部7により構成される特定話者音声認識部4と、音声信号入力部1と、映像信号入力部2と、不特定話者音声認識部3と、認識結果統合部5と、顔領域抽出部9と、顔画像データベース11により構成される。

【0044】この装置では、音声入力部1に入力した話者の音声は不特定話者認識部3と特定話者認識部4内の音声処理部6に入力する。音声処理部6の出力は、音声認識処理部7に入力し、音声認識処理部7の出力は、特定話者音声認識部4の出力として出力される。また、音声認識処理部7と音声データベース8は互いに接続されている。映像信号入力部2に入力した話者の顔を含む映像は、顔領域抽出部9に入力し、顔領域抽出部9の出力は画像比較部10に入力し、画像比較部10の出力は音声認識処理部7に接続される。また、画像比較部10と画像データベース11とは互いに接続されている。不特定話者音声認識部3、特定話者音声認識部4、画像比較部10の出力は、認識結果統合部5に入力し、認識結果統合部5からは認識結果が出力される。記録制御部2には、画像比較部10、不特定話者音声認識部3、特定話者音声認識部4の出力が接続され、記録制御部2の出力は音声データベース8と顔画像データベース11に接続される。

【0045】顔画像比較部10では、登録された複数の話者に対して、 i 番目の話者の顔画像データと入力映像に含まれ顔領域抽出部9から出力する顔画像データとの間の類似度 R_i を算出し出力する。顔画像データベース11に登録されたすべての話者 i に対して、画像比較部10の出力 R_i が、あらかじめ設定された閾値以下の場合、記録制御部19は、新しい識別コードを付与して、顔領域抽出部9の出力を顔画像データベース11に記録するよう制御する。

【0046】特定話者音声認識部4では、登録された複数の話者のうち画像比較部の出力 R_i があらかじめ設定された閾値 S_1 以上となる話者に対して i 番目の話者の音声データ j と、入力音声との類似度を $R'_{i,j}$ を算出して出力する。不特定話者認識部3では、音声データ j に対して複数の不特定話者の音声から共通する特徴を抽出して作成した標準パターンと入力音声との間の類似度 R''_j を算出し、出力する。音声データ j に対して、特定話者音声認識部4の出力 $R'_{i,j}$ が、あらかじめ設定した閾値 S_2 よりも小さく、かつ、不特定話者音声認識部

3の出力 R''_j があらかじめ設定した閾値 S_3 よりも大きい場合には、記録制御部19は、音声データ番号 j 、話者の名前、あるいは話者識別コードとともに入力音声を音声データベースに記録するよう制御する。

【0047】このように入力音声の話者が未登録の場合は話者の顔画像データと音声データを、話者は登録済みで音声データが未登録の場合は音声データを自動的に追加記録することが可能となる。なお、上記実施の形態1から5では、認識結果統合部5の入力として画像比較部10からの出力を要件としている（方式1）が、それを必須とはしない不特定話者音声認識部3及び音声認識処理部7の2つの出力を入力（方式2）としても一向に構わない。

【0048】方式2による方法では、主に以下の場合に有効である。すなわち、（1）一般に画像処理は負荷が大きいのでそれを軽減する対策として、音声認識処理部7の処理対象を画像比較部10で類似度の大きいものに絞ることにより、音声認識処理部7の処理負荷を少しでも軽減したい場合、（2）音声データベース8のデータが膨大な為、（画像処理の負荷の大小に関係なく）音声認識処理部の負荷を軽減したい場合、などである。

【0049】一方、方式1では上記実施の形態1から5の内容に加え、方式2と比較した場合、例えば以下のような場合に有効である。すなわち、（3）例えば登録された顔画像が正面である場合には、画像比較の対象として横顔等が入力されると、正面から捉えた顔（外観的特徴）でない為、その類似度（画像比較部10からの出力）の信頼性はやや低下する。そのような場合、方式2のように音声データベース8を前記類似度で絞り込むと音声認識処理部7の出力の信頼性を低下させる可能性があるため、認識結果統合部5での統合化処理の優先度として、不特定話者音声認識部3及び音声認識処理部7の出力（絞り込みを行わない出力）を優先しつつ画像比較部10の出力も有効活用したい場合、などである。

【0050】もちろん、上記方式1、方式2を（自動）切り替えするようにすれば、様々な利用形態に対応したより信頼性の高い認識結果（認識結果統合部5の出力）が得られることは言うまでもない。

【0051】

【発明の効果】以上の説明から明らかなように、本発明の音声認識装置は話者の特徴的外観として、例えば顔などを含む映像から話者の顔画像を抽出して、登録された話者の顔画像データベースと照合し、類似度を算出して、特定話者音声認識部、不特定話者音声認識部から出力する音声の類似度との統合的な類似度を算出して認識結果を出力することにより、複数の特定話者の入力に対して、信頼性の高い音声認識を行うことが可能となる。

【図面の簡単な説明】

【図1】本発明の第1の実施形態における音声認識装置の概略構成を示すブロック図

【図2】本発明の第2の実施形態における音声認識装置の概略構成を示すブロック図

【図3】本発明の第3の実施形態における音声認識装置の概略構成を示すブロック図

【図4】本発明の第4の実施形態における音声認識装置の概略構成を示すブロック図

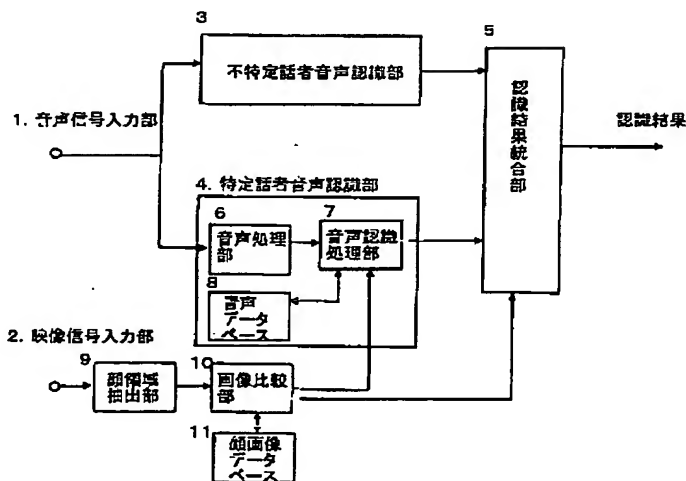
【図5】本発明の第5の実施形態における音声認識装置の概略構成を示すブロック図

【図6】従来の音声認識装置の概略構成図

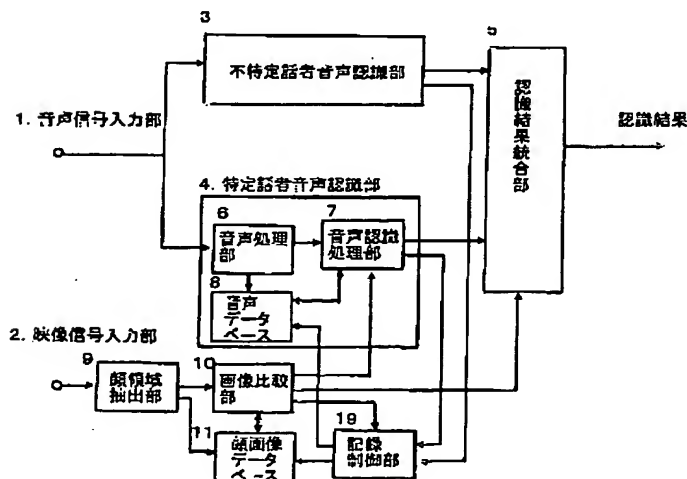
【符号の説明】

- | | |
|--------------|-------------|
| 1 音声信号入力部 | 6 音声処理部 |
| 2 映像信号入力部 | 7 音声認識処理部 |
| 3 不特定話者音声認識部 | 8 音声データベース |
| 4 特定話者音声認識部 | 9 顔領域抽出部 |
| 5 認識結果統合部 | 10 画像比較部 |
| | 11 画像データベース |
| | 12 唇領域抽出部 |
| | 13 唇動き検出部 |
| | 14 音声レベル検出部 |
| | 15 番組表入力部 |
| | 16 番組表記録部 |
| | 17 出演者名検出部 |
| | 18 画像検索部 |
| | 19 記録制御部 |

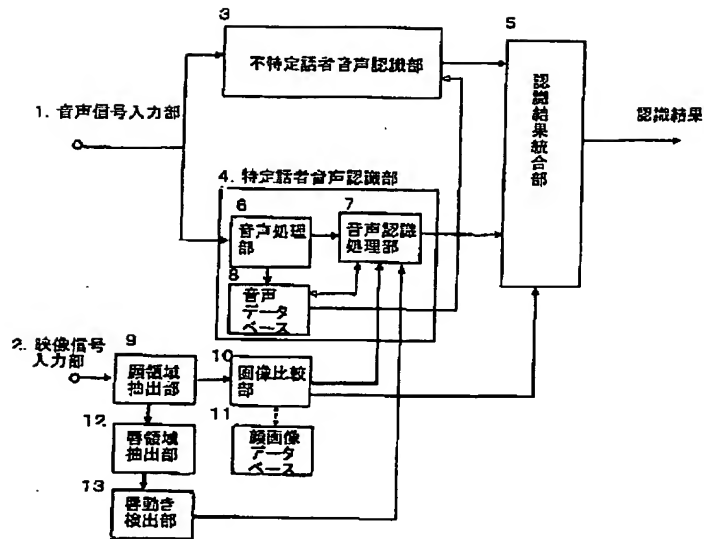
【図1】



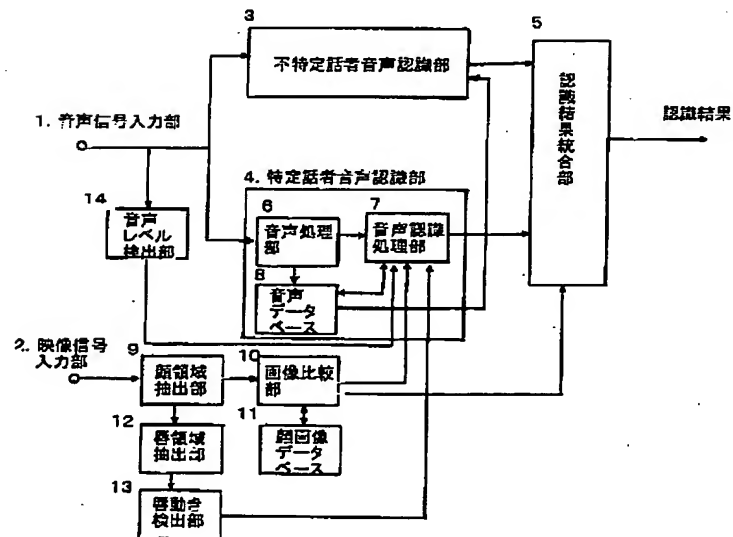
【図5】



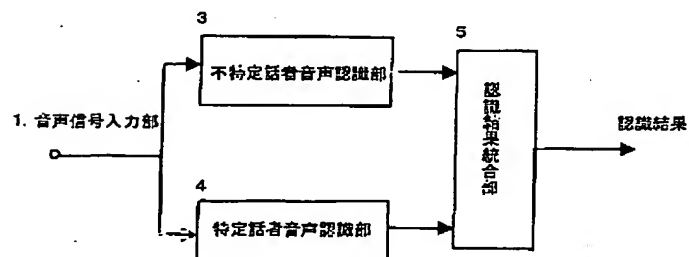
【図2】



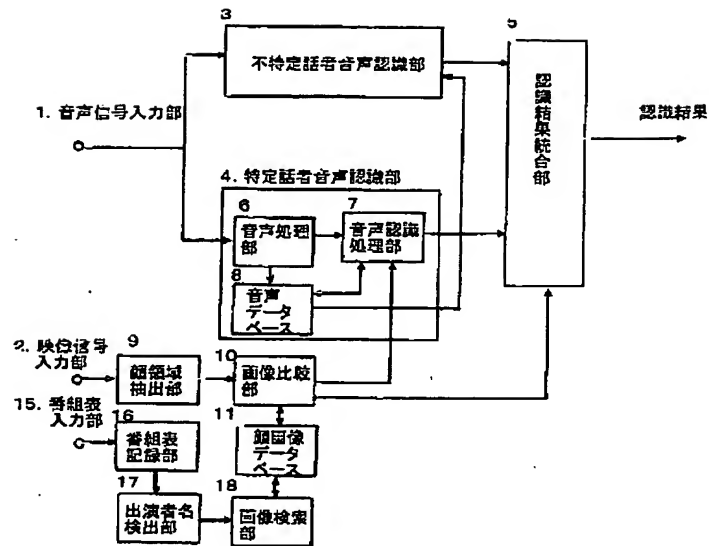
【図3】



【図6】



【図4】



フロントページの続き

(51)Int. Cl.⁶

G10L 5/06

識別記号

FI

G06F 15/62

380